

WHAT IS CLAIMED IS:

1. A method of node translation for communicating over virtual channels in a clustered multiprocessor system, the clustered multiprocessor system including a plurality of processing element nodes including a local processing element node and a remote processing element node, and a network interconnect coupled between the processing element nodes for sending communications between the processing element nodes, comprising:
 - assigning a connection descriptor to a virtual connection, the connection descriptor being a handle that specifies an endpoint node for the virtual connection;
 - defining a local connection table configured to be accessed using the connection descriptor to produce a system node identifier for the endpoint node;
 - generating a communication request including the connection descriptor;
 - in response to the communication request, accessing the local connection table using the connection descriptor of the communication request to produce the system node identifier for the endpoint node for the virtual connection; and
 - sending a memory request to the endpoint node, wherein the memory request is sent to the local processing element node if the endpoint node is the local processing element node, and is sent over the network interconnect to the remote processing element node if the endpoint node is the remote processing element node.
2. The method of claim 1, wherein assigning is performed by a local operating system in response to an operating system call by a local user process.
3. The method of claim 2, wherein assigning uses the connection descriptor to define a logical connection between a first virtual address space used by a local user process and a second virtual address space, whereby the connection descriptor allows the local user process to access the second virtual address space.

4. The method of claim 2, wherein defining is performed by the local operating system, which defines a local connection table for each user process.

5. The method of claim 1, wherein the local connection table includes a plurality of entries which are indexed by the connection descriptor, with each entry providing the system node identifier for the endpoint node of the associated virtual connection, and a key for qualifying an address translation at the endpoint node.

6. The method of claim 5, wherein each entry of the local connection table also includes a valid field for indicating whether the local connection table entry is valid.

7. The method of claim 1, wherein generating is performed by a local user process, which applies the communication request to a communication engine.

15 8. The method of claim 7, wherein accessing and sending are performed by the communication engine in response to the communication request.

9. The method of claim 5, wherein sending includes sending a virtual address and the key to the identified endpoint node, the address and key being for use in performing and qualifying address translation at the endpoint node.

20 10. The method of claim 1, wherein generating, accessing and sending are performed without intervention by an operating system.

25 11. The method of claim 1, wherein sending includes determining if a maximum number of outstanding packets has been reached.

12. A method of node translation for communicating over virtual channels in a clustered multiprocessor system, the clustered multiprocessor system including a plurality of processing element nodes, including a local processing element node

30

and a remote processing element node, and a network interconnect coupled between the processing element nodes for sending communications between the processing element nodes, comprising:

- 5 assigning a first connection descriptor to a first virtual connection having a first endpoint node which is a local node on the local processing element node;
- assigning a second connection descriptor to a second virtual connection having a second endpoint node which is a remote node on the remote processing element node;
- 10 defining a local connection table configured to be accessed using the first and the second connection descriptors to produce a first and a second system node identifier for the first and the second endpoint nodes, respectively;
- generating a communication request including a target connection descriptor;
- 15 in response to the communication request, accessing the local connection table to produce the first system node identifier if the target connection descriptor is the first connection descriptor, and to produce the second system node identifier if the target connection descriptor is the second connection descriptor; and
- 20 sending a memory request to the endpoint node identified by accessing the local connection table.

20 13. The method of claim 12, wherein assigning the first and the second connection descriptors are performed by a local operating system in response to an operating system call.

25 14. The method of claim 13, wherein assigning the second connection descriptor uses the second connection descriptor to define a logical connection between a first virtual address space used by a local user process and a second virtual address space, whereby the second connection descriptor allows the local user process to access the second virtual address space.

15. The method of claim 13, wherein defining is also performed by the local operating system, which defines a local connection table for each user process.

16. The method of claim 12, wherein the local connection table includes a plurality of entries indexed by the connection descriptors, with each entry providing the system node identifier for the endpoint node of the associated virtual connection, and also providing a key for qualifying an address translation at the endpoint node.

17. The method of claim 16, wherein each entry of the local connection table also includes a valid field for indicating if that local connection table entry is valid.

18. The method of claim 12, wherein generating is performed by a local user process, which applies the communication request to a communication engine.

15 19. The method of claim 18, wherein accessing and sending are performed by the communication engine in response to the communication request.

20. The method of claim 16, wherein sending includes sending a virtual address and the key to the endpoint node of the identified virtual connection for use in performing and qualifying an address translation at the endpoint node.

21. The method of claim 12, wherein generating, accessing and sending are performed without intervention by an operating system.

25 22. In a clustered multiprocessor system including a plurality of processing element nodes, including a local processing element node and a remote processing element node, and a network interconnect coupled between the processing element nodes for sending communications between the processing element nodes, a node translation apparatus comprising:

a memory configured to store a local connection table having a plurality of entries indexed by a connection descriptor, each entry of the local connection table storing a system node identifier for the endpoint of a virtual connection; and

5 a communication engine configured to receive a communication request including a connection descriptor from a user process, to access the local connection table using the connection descriptor of the communication request to produce the system node identifier for the endpoint node for the virtual connection, and to send a memory request to the endpoint node identified using the local connection table, wherein the memory request is sent internally to the endpoint node if the endpoint node is located within the local processing element node, and is sent over the network interconnect to the endpoint node if the endpoint node is located within the remote processing element node.

10

23. A clustered multiprocessor system, comprising:

15 a plurality of processing element nodes, including a local processing element node and a remote processing element node;

a network interconnect coupled between the processing element nodes for sending communications between the processing element nodes; and

20 a node translation apparatus having:

a memory configured to store a local connection table having a plurality of entries indexed by a connection descriptor, each entry of the local connection table storing a system node identifier for the endpoint of a virtual connection; and

25 a communication engine configured to receive a communication request including a connection descriptor from a user process, to access the local connection table using the connection descriptor of the communication request to produce the system node identifier for the endpoint node for the virtual connection, and to send a memory request to the endpoint node identified using the local connection table, wherein the memory request is sent internally to the endpoint node if the endpoint node is located within the

30

local processing element node, and is sent over the network interconnect to the endpoint node if the endpoint node is located within the remote processing element node.

5 24. The clustered multiprocessor system of claim 23, wherein the network interconnect includes two request and two response virtual channels.

25. The clustered multiprocessor system of claim 23, wherein the remote node includes means for suppressing TLB miss servicing as a function of aging of
10 outstanding packets.

26. The clustered multiprocessor system of claim 23, wherein each processing element node includes a communication engine, wherein address translation occurs within the communication engine.

15 27. The clustered multiprocessor system of claim 23, wherein a communication engine on one processing element node performs address translation for an address on another processing element node.